

PATENT

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re application of

Yuji KOBAYASHI

Serial No.: 10/762,126

Filed: January 21, 2004



Group Art Unit:

Examiner:

For: INFORMATION SEARCHING APPARATUS AND METHOD, INFORMATION
SEARCHING PROGRAM, AND STORAGE MEDIUM STORING THE
INFORMATION SEARCHING PROGRAM

Certificate of Mailing

I hereby certify that this paper is being deposited with the
United States Postal Service as first class mail in an
envelope addressed to: Commissioner for Patents, P.O.
Box 1450, Alexandria, VA 22313-1450 on:

Date: 02/05/04

By: [Signature]
Marc A. Rossi

CLAIM FOR PRIORITY

Commissioner for Patents
P.O. Box 1450
Alexandria, VA 22313-1450

Sir:

The benefit of the filing date of the following prior foreign application filed in the
following country is hereby requested for the above-identified application and the priority
provided in 35 U.S.C. § 119 is hereby claimed:

JAPAN 2003 - 013428 January 22, 2003

In support of this claim, a certified copy of said original foreign application is filed
herewith. It is requested that the file of this application be marked to indicate that the
requirements of 35 U.S.C. 119 have been fulfilled and that the Patent and Trademark Office
kindly acknowledge receipt of this document.

02/05/04
Date

Attorney Docket: CANO:116

Respectfully submitted,

[Signature]
Marc A. Rossi
Registration No. 31,923

日 本 国 特 許 庁
JAPAN PATENT OFFICE

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日 2 0 0 3 年 1 月 2 2 日
Date of Application:

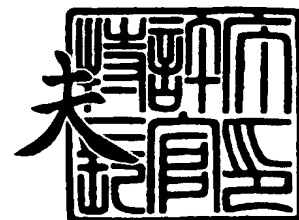
出 願 番 号 特 願 2 0 0 3 - 0 1 3 4 2 8
Application Number:
[ST. 10/C]: [J P 2 0 0 3 - 0 1 3 4 2 8]

出 願 人 キヤノン株式会社
Applicant(s):

2 0 0 4 年 1 月 1 4 日

特許庁長官
Commissioner,
Japan Patent Office

今 井 康 夫



【書類名】 特許願

【整理番号】 224320

【提出日】 平成15年 1月22日

【あて先】 特許庁長官殿

【国際特許分類】 G06F 1/00

【発明の名称】 情報検索装置

【請求項の数】 2

【発明者】

 【住所又は居所】 東京都大田区下丸子 3 丁目 3 0 番 2 号 キヤノン株式会社
社内

 【氏名】 小林 雄二

【特許出願人】

 【識別番号】 000001007

 【氏名又は名称】 キヤノン株式会社

 【代表者】 御手洗 富士夫

【代理人】

 【識別番号】 100081880

 【弁理士】

 【氏名又は名称】 渡部 敏彦

 【電話番号】 03(3580)8464

【手数料の表示】

 【予納台帳番号】 007065

 【納付金額】 21,000円

【提出物件の目録】

 【物件名】 明細書 1

 【物件名】 図面 1

 【物件名】 要約書 1

 【包括委任状番号】 9703713

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 情報検索装置

【特許請求の範囲】

【請求項 1】 文書から抽出された単語を前記文書と対応付けて文書検索用の索引情報として登録する索引情報登録手段と、該索引情報登録手段により登録された索引情報を参照して、要求された検索用情報に対応する文書を検索する文書検索手段とを有する情報検索装置であって、

前記文書から未知単語を抽出する未知単語抽出手段と、

前記文書の種別を判別する文書種別判別手段と、

前記文書種別判別手段により判別された文書の種別に応じて、前記未知単語抽出手段により抽出された未知単語の前記索引情報登録手段による索引情報としての登録の可否を決定する登録可否決定手段とを有することを特徴とする情報検索装置。

【請求項 2】 文書から単語を抽出する単語抽出手段と、

前記単語抽出手段により抽出された単語を前記文書と対応付けて、文書検索用の索引情報として登録する索引情報登録手段と、

前記索引情報登録手段により登録された索引情報を参照して、要求された検索用情報に対応する文書を検索する文書検索手段と、

文字認識処理された文書に含まれている文字コードの誤りを校正する文字校正手段とを有し、

前記単語抽出手段は、前記文字校正手段により校正された後の文書から単語を抽出することを特徴とする情報検索装置。

【発明の詳細な説明】

【0 0 0 1】

【発明の属する技術の分野】

本発明は、文字、画像等のマルチメディア情報を含んだ文書から所望の文書を検索する情報検索装置に関する。

【0 0 0 2】

【従来の技術】

従来、マルチメディア情報を含んだ文書から所望の文書を検索するため全文検索と呼ばれる手法を適用した情報検索装置が知られている。この装置では、例えば、所望の検索用情報を表す検索語あるいは文等を入力し、蓄積された文書群から、入力された語あるいは文と一致する語あるいは文を有している文書を得るようになっている。

【0 0 0 3】

また、文書画像であっても文書の内容に基づく情報検索を可能にするために、文書画像に含まれている文字画像部分を文字認識し、文字認識結果得られる文字情報をもとに全文検索を適用可能なように構成した装置が考案されている。

【0 0 0 4】

しかし、文字認識結果得られた文字コードを含んだ文書（以下、「文字認識処理済み文書」と称する）には、誤認識された文字が含まれる可能性があるため、文字認識を経ないテキスト等の文書と同様の手法によって全文検索を行なった場合、検索漏れの増大や、本来の文書内容と異なる文字との思いがけない一致による誤検索が生じることがある。

【0 0 0 5】

そのため、文字認識処理済み文書に全文検索を適用する前に、検索対象となる文字認識処理済み文書の誤認識箇所をひとつひとつユーザが目視して修正することが通常行われる。

【0 0 0 6】

そしてこの目視修正を不要とするために、下記特許文献1では、文字認識の候補となる複数の文字をその確からしさの推定値とともに用いることで、誤認識が含まれる文字認識処理済み文書であっても、複数の候補文字を選択して検索漏れを低減させる手法が開示されている。すなわち、複数の文字認識候補文字を含めて検索することで検索漏れを低減させることができる。

【0 0 0 7】

【特許文献1】

特許第 2 5 8 6 3 7 2 号

【0 0 0 8】

【発明が解決しようとする課題】

しかしながら、例えば、図 5 に示される誤認識文字列「モルール」に対して、たまたま誤認識文字列と一致してしまう「ルール」と照合することを回避できないため、文字列「ルール」で検索すると、本来は「モノレール」であるこの誤認識文字列「モルール」と一致し、誤検索されてしまう等のように、精度の低下を招く恐れがある。

【0 0 0 9】

また、文字認識により得られた文字コードのみで構成される文字認識処理済み文書の場合は、上記特許文献 1 の手法を用いたとしても、文字認識過程における他の候補文字の情報が必要であるため、その効果を期待できず、検索漏れの増大や誤検索の問題が残る。

【0 0 1 0】

一方で、単に文字単位の索引情報のみならず、実際に存在する単語との照合を行う形態素解析を行って、抽出された単語を文書検索用の索引情報として登録する、単語索引による情報検索装置も提案されている。この単語索引による情報検索装置では、文字単位の索引による情報検索装置と比べ、単語の境界をまたいだ一致等を避けることができるため、検索精度を向上させることができる。ところが、この単語索引による情報検索においても、すべての単語を単語辞書に収録することは事実上不可能であるため、辞書に存在しない単語について検索することができず、検索漏れが起こり得る。

【0 0 1 1】

本発明は上記従来技術の問題を解決するためになされたものであり、その第 1 の目的は、検索漏れ及び誤検索の少ない適切な検索を行うことができる情報検索装置を提供することにある。

【0 0 1 2】

本発明の第 2 の目的は、検索精度を向上させることができる情報検索装置を提供することにある。

【0 0 1 3】**【課題を解決するための手段】**

上記第 1 の目的を達成するために本発明の請求項 1 の情報検索装置は、文書から抽出された単語を前記文書と対応付けて文書検索用の索引情報として登録する索引情報登録手段と、該索引情報登録手段により登録された索引情報を参照して、要求された検索用情報に対応する文書を検索する文書検索手段とを有する情報検索装置であって、前記文書から未知単語を抽出する未知単語抽出手段と、前記文書の種別を判別する文書種別判別手段と、前記文書種別判別手段により判別された文書の種別に応じて、前記未知単語抽出手段により抽出された未知単語の前記索引情報登録手段による索引情報としての登録の可否を決定する登録可否決定手段とを有することを特徴とする。

【 0 0 1 4 】

上記第 2 の目的を達成するために本発明の請求項 2 の情報検索装置は、文書から単語を抽出する単語抽出手段と、前記単語抽出手段により抽出された単語を前記文書と対応付けて、文書検索用の索引情報として登録する索引情報登録手段と、前記索引情報登録手段により登録された索引情報を参照して、要求された検索用情報に対応する文書を検索する文書検索手段と、文字認識処理された文書に含まれている文字コードの誤りを校正する文字校正手段とを有し、前記単語抽出手段は、前記文字校正手段により校正された後の文書から単語を抽出することを特徴とする。

【 0 0 1 5 】

【発明の実施の形態】

以下、本発明の実施の形態を図面を参照して説明する。

【 0 0 1 6 】

図 1 は、本発明の一実施の形態に係る情報検索装置の全体構成を示すブロック図である。

【 0 0 1 7 】

図 1 において、11 はマイクロプロセッサ（CPU）（索引情報登録手段、文書検索手段、未知単語抽出手段、単語抽出手段、文書種別判別手段、文字認識処理手段、登録可否決定手段、文字校正手段）であり、情報検索のための演算、論理判断等を行い、アドレスバス AB、コントロールバス CB、データバス DB を

介して、それらのバスに接続された各構成要素を制御する。

【0 0 1 8】

アドレスバス A B は C P U 1 1 が制御の対象とする構成要素を指示するアドレス信号を転送する。コントロールバス C B は C P U 1 1 が制御の対象とする各構成要素のコントロール信号を転送して印加する。データバス D B は各構成要素相互間のデータ転送を行う。

【0 0 1 9】

1 2 は読出し専用の固定メモリ（R O M）であり、本実施の形態で実行される処理プログラム等の制御プログラムコードを記憶する。1 3 は書込み可能なランダムアクセスメモリ（R A M）であり、各構成要素からの各種データの一時記憶に用いられる。R A M 1 3 はまた、図 2 において後述する検索語保持部 2 0 2、検索結果保持部 2 0 4、未知語検索指定保持部 2 0 8（可否指定手段）、文書タイプ保持部 2 1 4 を記憶する。

【0 0 2 0】

1 4 は外部メモリ（D I S K）であり、図 2 において後述する単語インデックス 2 0 5、単語抽出辞書 2 0 7（単語辞書）、文字校正用辞書 2 1 0、抽出テキスト 2 1 2、文書登録処理部 2 1 5 に入力された登録文書が格納される。また、図 2 において後述する検索要求入力処理部 2 0 1、検索処理部 2 0 3（文書検索手段）、単語抽出処理部 2 0 6（未知単語抽出手段、単語抽出手段）、文字校正処理部 2 0 9（文字校正手段）、文字認識処理部 2 1 1（文字認識処理手段）、登録文書タイプ判定部 2 1 3（文書種別判別手段）の各処理部を実行するプログラムコードが格納される。また、これらのデータ及びプログラムを格納する記憶媒体としては、固定磁気ディスク、R O M、フロッピー（登録商標）ディスク、C D - R O M、C D - R、C D - R W、D V D - R O M、D V D - R A M、D V D - R、メモ리카ード、光磁気ディスク、磁気テープ等を用いることができる。

【0 0 2 1】

1 5 はキーボード（K B）であり、アルファベットキー、ひらがなキー、カタカナキー、句点等の文字記号入力キー、検索を指示する検索キー及び、カーソル移動を指示するカーソル移動キー等のような各種の機能キーを備えている。1 6

は表示用ビデオメモリ（VRAM）であり、表示すべきデータのパターンを蓄える。17はCRTコントローラ（CRTC）であり、表示用ビデオメモリVRAM16に蓄えられた内容を後述する表示器CRT18に表示する役割を担う。

【0022】

表示器CRT18は、陰極線管、液晶パネル等の表示装置であり、表示装置CRT18におけるドット構成の表示パターン及びカーソルの表示は、CRTコントローラ17により制御される。19はネットワークコントローラ（NIC）であり、Ethernet（登録商標）等のネットワークに接続する役割を担う。

【0023】

かかる各構成要素からなる情報検索装置は、キーボードKB15からの各種の入力及びネットワークコントローラ19から供給されるネットワーク経由の各種入力に応じて作動し、キーボードKB15からの入力及びネットワークコントローラ19からの入力供給されると、まず、インタラプト信号がCPU11に送られ、CPU11がDISK14内に記憶されている各種の制御信号を読み出し、それらの制御信号に従って、各種の制御を行う。

【0024】

図2は、本発明の実施の形態の情報検索装置の機能構成を示すブロック図である。

【0025】

同図において、201は所望の検索対象に関する要求事項（検索文や検索語等の検索用情報であって、本実施の形態では、検索語を用いることとする）を入力する検索要求入力処理部である。202は検索要求入力処理部201によって入力された検索語を記憶する検索語保持部である。203は検索語保持部202に記憶された検索語に基づいて、登録された文書を対象として検索する検索処理部である。204は、検索処理部203の処理結果を保持する検索結果保持部である。

【0026】

205は、文書登録処理部215で登録された文書から抽出された単語と、その単語の出現する文書の文書情報とを対応させて記憶した単語インデックスであ

る（後述）。2 0 6 は、文書登録処理部 2 1 5 で登録された文書から検索対象とする単語（検索語との照合対象となる単語）を抽出する単語抽出処理部である。2 0 7 は、単語抽出処理部 2 0 6 で参照される単語を定義した単語抽出辞書である。

【 0 0 2 7 】

2 0 8 は、文書登録処理部 2 1 5 で登録された文書において出現した単語のうち単語抽出辞書 2 0 7 において定義されていない「未知語」とであると判定された単語に対して、これを単語インデックス 2 0 5 に、文書検索用の索引情報として登録する否かを指定する未知語検索指定保持部である。この未知語検索指定保持部 2 0 8 では、文字認識処理を経た結果であるテキスト文書（すなわち、文字認識の結果得られた文字コードを含んだ文書であり、以下、「文字認識処理済み文書」とも称する）の登録の際に、単語抽出辞書 2 0 7 に収録されていない未知語であっても検索可能なように登録する（インデックスを作成する）か、あるいは未知語については検索対象とならないように登録を禁止するかが、ユーザの意思で指定される。すなわち、未知語検索指定保持部 2 0 8 は、文字認識処理済み文書から抽出された未知語を索引情報として登録することを許可するか否かの指定を記憶するものである。

【 0 0 2 8 】

2 0 9 は、文字認識処理部 2 1 1 において文書画像から得られた認識文字よりなる抽出テキスト 2 1 2 に対して、文字校正用辞書 2 1 0 及び単語抽出辞書 2 0 7 を参照して文字認識の誤りを校正する文字校正処理部である。2 1 1 は、文書登録処理部 2 1 5 で登録された文書が画像文書である場合に、その画像情報を文字認識処理により文字情報に変換する文字認識処理部である。2 1 2 は、文字認識処理部 2 1 1 の処理結果である認識文字（文字認識処理済み文書を構成する文字コード）を格納する抽出テキストである。

【 0 0 2 9 】

2 1 3 は、文書登録処理部 2 1 5 で登録された文書の文書タイプ（画像文書であるか、文字認識結果を格納した抽出テキストであるか、一般テキストであるか）を判定する登録文書タイプ判定部である。登録文書タイプ判定部 2 1 3 では、

与えられた登録文書のファイル名の拡張子が”bmp”、”jpg”、”gif”、”tif”等の画像フォーマットをあらわす場合には、文書タイプ保持部 2 1 4 に「画像タイプ」であることを格納し、登録文書のファイル名の拡張子が”txt”等、テキストをあらわす場合には、文書タイプ保持部 2 1 4 に「テキストタイプ」であることを格納する。また、テキスト文書でも特に文字認識処理の結果得られた抽出テキストである場合には、特別な拡張子として”ocr”を使用するものとし、拡張子が”ocr”であった場合は、文書タイプ保持部 2 1 4 に「認識結果テキストタイプ」であることを格納する。2 1 5 は文書の登録を行う文書登録処理部である。

【 0 0 3 0 】

次に、検索要求入力処理部 2 0 1 において、文書検索の検索要求のための検索語を操作者が指示する場合における操作パネルの表示例について、図 3 を用いて説明する。

【 0 0 3 1 】

図 3 は、検索語を指示する場合において、表示器 C R T 1 8 に画面表示される操作パネルの例を示す図である。

【 0 0 3 2 】

図 3 において、3 0 1 は検索要求入力操作を行う表示ウインドウである。3 0 2 は検索語等の検索用情報を入力する検索文入力領域である。3 0 3 は入力中の検索用情報（検索要求文）の一例を示しており、同図においては、「カラープリンタの売上」と入力されている。3 0 4 は検索文入力領域 3 0 2 における入力位置を示す入力カーソルである。

【 0 0 3 3 】

3 0 8 は検索処理の実行を指定する検索実行ボタンであり、検索実行ボタン 3 0 8 を押下することで、指定した検索処理が実行される。3 0 9 は検索処理の終了あるいは中止を指定するキャンセルボタンであり、キャンセルボタン 3 0 9 を押下すると、ただちに検索処理が終了し、表示ウインドウ 3 0 1 が閉じられて終了する。3 1 0 は検索ボタン 3 0 8 の押下によって検索処理を行った結果を表示する検索結果表示領域であり、同図においては検索処理がなされていない状態であるので、何も表示されていない。

【 0 0 3 4 】

図 4 は、検索要求入力処理部 2 0 1 による検索処理実行後の検索結果の表示例を示す図である。同図には、検索用情報 3 0 3 に示された「カラープリンタの売上」に対して、文書検索を行った検索結果が表示されている。

【 0 0 3 5 】

図 4 において、4 0 2 は、検索結果の順位を示すランク表示領域である。検索結果は検索要求に適合している順にランク付けされ、ランク順に表示される。図 4 の表示例においては、ランク 2 5 位から 3 0 位までの検索結果が表示されている。4 0 3 は、検索された文書の表題であり、4 0 4 は文書のファイル名である。4 0 5 は、検索された文書の大意が掴める程度の内容を表示する文書内容表示である。文書内容表示 4 0 5 には、あらかじめ文書の書誌的属性として与えられた要約文、あるいは文書内容を自動的に要約した要約文、あるいは文書の一部を大意として抽出した大意文等を表示することができる。

【 0 0 3 6 】

4 0 6 は、表示位置を指定するために同種のウインドウ表示装置において用いられているエレベータバーであり、検索結果表示領域 3 1 0 に表示しきれない場合に、検索結果表示領域 3 1 0 内において検索結果を部分表示しながら、表示されていない他の部分を表示するために用いられる。

【 0 0 3 7 】

次に、文字認識処理を施した認識結果である抽出テキスト 2 1 2 について、図 5 を用いて説明する。

【 0 0 3 8 】

図 5 は、文字を表す画像に対して文字認識処理を施した結果得られた抽出テキストを示す図である。抽出テキストは、文字認識処理部 2 1 1 において抽出される場合と、あらかじめ文字認識処理が施されて文書登録処理部 2 1 5 に与えられている場合がありえる。

【 0 0 3 9 】

一般に、文字認識処理においては、認識誤りを含むことがありえる。例えば、同図に示すように、「モノクロ」の 4 文字目「ロ」がカタカナの「ロ」（ろ）で

あるべきところ、漢字の「口」（くち）と誤認識され、文字認識処理の対象とされた元画像では「モノレール」と書かれていた部分が「モルール」と誤認識されている。

【0 0 4 0】

次に、文字認識結果である抽出テキストに対して、文字校正処理を行った処理結果の例について図 6 を用いて説明する。

【0 0 4 1】

図 6 は、図 5 に例示された文字認識誤りを含んでいる抽出テキストに対して文字校正処理を行った場合の結果を示す図である。文字校正処理部 2 0 9 では、後述する文字校正用辞書 2 1 0 と単語抽出辞書 2 0 7 を参照することにより、図 5 における誤認識である「モノクロ」（「口」は漢字の「くち」）を「モノクロ」（「口」はカタカナ）と校正している。

【0 0 4 2】

図 7 は、文字校正用辞書 2 1 0 の構成を示す概念図である。文字校正用辞書 2 1 0 は、文字認識処理において誤認識の発生しやすい、似通った字形の文字を対応付けて格納したものである。同図において、1 つの行に表されている文字同士が、互いに誤認識されやすいことを表している。

【0 0 4 3】

図 8 は、単語抽出辞書 2 0 7 の構成を示す概念図である。同図において、8 0 1 は単語の見出し語、8 0 2 は単語の品詞を示す。品詞 8 0 2 は、抽出された単語と前置される単語との接続可能性を判定するために用いられる。

【0 0 4 4】

図 9 は、単語インデックス 2 0 5 の構成を示す概念図である。単語インデックス 2 0 5 は、登録文書中に出現するすべての単語について、文書中における出現頻度とその文書との対応をとって格納するテーブルである。単語インデックス 2 0 5 におけるテーブルの第 1 列情報である 9 0 1 は、検索見出しとなる単語（インデックス見出し語）である。テーブルの第 2 列情報は、インデックス見出し語 9 0 1 が出現する文書と、その文書においてインデックス見出し語 9 0 1 が何回出現するかを対応づけた文書情報であり、インデックス見出し語 9 0 1 が出現す

るすべての文書について格納される。その際、文書は一意化された文書識別番号で記録される。

【0 0 4 5】

例えば、インデックス見出し語 9 0 1 「カラー」に対応して、文書情報 9 0 2 において (1 0 0 0、1 5)、(1 2 0 0、5) 等と記録されており、これによれば、「カラー」は、文書識別番号「1 0 0 0」で表される文書に 1 5 回出現し、文書識別番号「1 2 0 0」で表される文書に 5 回出現していることが示される。

【0 0 4 6】

図 1 0 は、文字校正処理で用いられる文字候補ラティスの構成を示す概念図である。

【0 0 4 7】

文字候補ラティスは、文字校正処理部 2 0 9 での文字校正処理において、誤認識の発生する可能性のある文字について、他の候補文字を格子状に配したものである。文字校正処理対象の抽出テキストの各文字について、図 7 に示す文字校正用辞書 2 1 0 に格納された文字のいずれかと一致したならば、対応している他の文字が置換候補文字として追加される。

【0 0 4 8】

図 1 0 においては、図 5 に示す文字認識結果の抽出テキストに対して作成された文字候補ラティスの一部が例示されており、「このモノクロの写真に写っているモルルは当時の重要な交通手段であった」という文字列に対して、図 7 の文字校正用辞書 2 1 0 に出現する「口」と「一」に対して、候補文字を追加したラティスが作成されている。

【0 0 4 9】

図 1 1 は、本実施の形態における文書検索処理のフローチャートを示す図である。

【0 0 5 0】

まず、ステップ S 1 0 0 1 では、検索要求入力処理部 2 0 1 の動作を行う処理モジュールによって、検索要求入力処理を行う。検索要求入力処理では、図 3 に

示す操作パネル中の検索文入力領域 302 に入力された検索用情報（本実施の形態では検索語）を取り出し、単語抽出辞書 207 を参照して、入力検索用情報から単語を抽出し、抽出した検索語を検索語保持部 202 に格納する。

【0051】

次に、ステップ S1002 では、検索語保持部 202 に格納された検索語が含まれる文書を単語インデックス 205 を参照して検索する。すなわち、検索語保持部 202 より検索語を取り出し、取り出された検索語と一致するインデックス見出し語 901 を検索し、見つかった見出し語に対応する文書情報 902 を取り出す。そして、取り出された文書情報 902 中の文書識別番号と出現頻度とを検索結果保持部 204 に格納する。なお、既に同一の文書識別番号をもつ文書情報が検索結果保持部 204 に格納済みであった場合は、その出現頻度を更新する。検索語保持部 202 に格納されたすべての検索語についてこの処理を行い、処理が終わったならば、検索語保持部 202 に記憶された検索結果を出現頻度の大きい順にソートする。

【0052】

次に、ステップ S1003 では、前記ステップ S1002 で検索された検索結果を、検索結果保持部 204 より取り出して表示する。なお、この処理は同種の情報検索装置において広く行われている公知の処理と同様になされる。その後、本処理を終了する。

【0053】

図 12 は、本実施の形態における文書登録処理のフローチャートを示す図である。

【0054】

まず、ステップ S3001 では、文書登録処理部 215 に入力された文書の文書タイプを判定する。文書タイプの判定は、文書登録処理部 215 に入力された文書ファイル名の拡張子で行い、拡張子が "bmp"、"jpg"、"gif"、"tif" 等である場合は、「文書画像」であるので、文書タイプ保持部 214 に「画像タイプ」であることを格納するとともに、ステップ S3002 へ進む。また、拡張子が "ocr" であれば、文字認識処理の結果得られた抽出テキストであるので、文書タ

イブ保持部 2 1 4 に「文字認識抽出テキストタイプ」であることを格納すると共に、抽出テキスト 2 1 2 にその内容を格納して、ステップ S 3 0 0 3 へ進む。一方、拡張子が“txt”や“html”等である場合は、一般テキストであるので、文書タイプ保持部 2 1 4 に「テキストタイプ」であることを格納するとともに、ステップ S 3 0 0 4 へ進む。

【 0 0 5 5 】

ステップ S 3 0 0 2 では、画像タイプと判定された文書に対して、その画像中の文字部分について文字認識処理を行う。この場合、文書タイプが画像タイプから文字認識抽出テキストタイプに変わったので、認識文字の抽出テキストを抽出テキスト 2 1 2 に作成し、文書タイプ保持部 2 1 4 に「文字認識抽出テキストタイプ」であることを格納して、ステップ S 3 0 0 3 に進む。なお、画像情報から文字画像と照合を行い、文字コード化を行う文字認識処理は公知の手法でなされる。

【 0 0 5 6 】

ステップ S 3 0 0 3 では、抽出テキスト 2 1 2 に対して、文字認識の誤り訂正を行うために、後述する図 1 3 の文字校正処理を行う。ここでは、前記ステップ S 3 0 0 2 の文字認識処理の結果得られた認識文字抽出テキスト、あるいは前記ステップ S 3 0 0 1 において、「文字認識抽出テキストタイプ」と判定された文書である抽出テキストに対して文字校正処理がなされる。

【 0 0 5 7 】

続くステップ S 3 0 0 4 では、後述する図 1 4 の単語抽出処理を実行する。すなわち、文字認識結果ではない一般の文書（一般テキスト）、及び前記ステップ S 3 0 0 3 において文字校正処理の施された文書（文字認識抽出テキスト）から、単語を抽出し、文書検索のための単語インデックス 2 0 6 を作成する。

【 0 0 5 8 】

図 1 3 は、図 1 2 のステップ S 3 0 0 3 で実行される文字校正処理のフローチャートを示す図である。

【 0 0 5 9 】

まず、ステップ S 2 0 0 1 では、文書校正処理で用いる文字候補ラティスを作

成する。文字候補ラティスは、前述のように、文字認識結果の抽出テキストのうち、誤認識の起こりやすい文字について候補文字を加えて、ラティス上に構成したものである（図 1 0 参照）。文字候補ラティス作成処理では、対象となる抽出テキスト 2 1 2 より文字を取り出し、文字校正用辞書 2 1 0 を参照し、取り出した文字が文字校正用辞書 2 1 0 に登録されている文字であった場合は、その文字を含む、似た字形のグループの各文字を候補文字として文字候補ラティスに追加する。その際、ラティスの各格子点の候補文字の 1 番目の文字としては元の抽出テキストの文字を配するものとする。

【 0 0 6 0 】

次に、ステップ S 2 0 0 2 では、文字候補ラティスの各格子点の 1 番目の文字に従って単語抽出辞書 2 0 7 を検索する。すなわち、ラティスを構成する各格子に位置する 1 番目の文字列と一致する見出し語があるかどうかを検索する。

【 0 0 6 1 】

次に、ステップ S 2 0 0 3 では、前記ステップ S 2 0 0 2 での単語検索の結果、単語が検索されたかどうか判別する。単語が検索された場合は、ステップ S 2 0 0 8 へ進む一方、単語が検索されなかった場合は、ステップ S 2 0 0 4 へ進む。

【 0 0 6 2 】

ステップ S 2 0 0 4 では、文字候補ラティスにまだ単語検索していない候補文字列があるかを判別し、未検索の候補文字列がなかったならば前記ステップ S 2 0 0 8 へ進む一方、未検索の候補文字列があったならばステップ S 2 0 0 5 へ進む。

【 0 0 6 3 】

ステップ S 2 0 0 5 では、文字候補を変更して、未検索の候補文字列の単語検索を行う。次に、ステップ S 2 0 0 6 では、前記ステップ S 2 0 0 5 での単語検索の結果、一致する単語が検索されたか否かを判別し、単語が検索されなかったならば、再び前記ステップ S 2 0 0 4 へ戻って、未検索の候補文字列があるかの判別を繰り返す。

【 0 0 6 4 】

一方、単語が検索されたならばステップ S 2007 へ進んで、前記ステップ S 2005 において検索された単語と一致する候補文字列を構成する文字に、抽出テキストの文字を置き換える。例えば、図 10 の文字候補ラティスに対して、図 8 の単語抽出辞書 207 を検索した場合、候補文字列「モノクロ」と一致するため、抽出テキスト中の「モノクロ」の「ロ」（漢字の「ロ」）を「ロ」（カタカナの「ロ」）に置換する。

【0065】

次に、ステップ S 2008 では、次の単語検索位置を取得する。すなわち、単語抽出辞書 207 との照合の終わった文字候補列をスキップし、未照合の候補文字列開始位置を取得する。ここで、前記ステップ S 2004 から分岐してきた場合は、文字候補ラティス中の候補文字列と単語との照合がとれなかった場合であるので、最後に照合のとれた文字位置以降の最初の助詞の次の位置を、次の候補文字列開始位置とする。

【0066】

次に、ステップ S 2009 では、前記ステップ S 2008 で得られた候補文字列開始位置が文書末尾に達したか、すなわち、抽出テキストのすべての文字列について単語抽出辞書 207 との照合による文字校正処理を終えたかどうかを判別し、文書末尾に達していなければ前記ステップ S 2002 へ戻って未処理の文字列について処理を繰り返す一方、文書末尾に達していたならば本処理を終了する。

【0067】

図 14 は、図 12 のステップ S 3004 で実行される単語抽出処理のフローチャートを示す図である。

【0068】

まず、ステップ S 4001 では、単語検索、すなわち、単語抽出処理対象のテキストの文字列と、単語抽出辞書 207 との照合を行う。そして、ステップ S 4002 では、単語が検索されたか否かを判別する。その判別の結果、単語が検索されなかった場合は、ステップ S 4004 に進む一方、単語が検索された場合は、ステップ S 4003 へ進む。

【0069】

ステップS4003では、検索された単語と、既に抽出済みの直前の単語（前接語）が接続可能か否かを、単語抽出辞書207の品詞情報に基づき接続判定表（不図示）を用いて判別する。前接語との接続可能性の判定手法、及び接続可能性判定に用いる接続判定表の構成については公知であるため、詳細な説明を省略する。

【0070】

その判別の結果、抽出済みの単語と検索された単語とが接続可能である場合は、ステップS4007へ進む一方、接続不可能である場合は、ステップS4004へ進む。

【0071】

ステップS4004では、前記ステップS4001での単語検索の結果、単語が検索されなかったか、あるいは、前接語と接続可能な単語が検索されなかった場合であるので、その検索の照合開始位置から、辞書登録されていない未知語の抽出処理を行う。

【0072】

未知語の抽出処理は、例えば、連続するカタカナをひとつの未知語として抽出する、あるいは、照合開始位置の文字から頻度の高い助詞の出現する直前の文字までをひとつの未知語として抽出する等の公知の手法を適用可能である。この未知語抽出手法を、図6に示される文字校正処理を行ったテキストに対して適用することで、図8に示される単語抽出辞書207には未登録である「モルール」が未知語として抽出される。なお、上述のように、図6における「モルール」は、本来は「モノレール」であるべきところ、誤認識されたものである。

【0073】

次に、ステップS4005では、単語抽出処理対象文書が文字認識結果文書、すなわち、文字認識処理済み文書であるか否かを判定する。文字認識結果文書であるかどうかは文書タイプ保持部214に記憶された文書タイプを参照することにより行われる。例えば、図12のステップS3001またはステップS3002において、文書タイプ保持部214に「文字認識抽出テキストタイプ」と記憶

されていれば、文字認識結果文書であると判別される。

【0 0 7 4】

その判別の結果、単語抽出処理対象文書が文字認識結果文書（文字認識処理済み文書）でない場合は、ステップ S 4 0 0 7 に進む一方、単語抽出処理対象文書が文字認識結果文書（文字認識処理済み文書）である場合は、ステップ S 4 0 0 6 へ進み、文字認識結果文書に対する未知語検索指定の有無を判定する。未知語検索指定の有無の判別は、未知語検索指定保持部 2 0 8 を参照することにより行われる。ここで、「未知語検索指定有り」は、文字認識処理済み文書から抽出された未知語を単語インデックス 2 0 5 に索引情報として登録することを許可する指定であり、「未知語検索指定無し」は禁止する指定である。

【0 0 7 5】

その判別の結果、「未知語検索指定有り」である場合は、前記ステップ S 4 0 0 7 へ進む。

【0 0 7 6】

ステップ S 4 0 0 7 では、前記ステップ S 4 0 0 1 で検索された（既知の）単語、あるいは前記ステップ S 4 0 0 4 で抽出された未知語を、単語インデックス 2 0 5 に索引情報として登録する。その際、登録しようとする単語が単語インデックス 2 0 5 のインデックス見出し語 9 0 1 に既に存在している場合において、インデックス見出し語 9 0 1 に対応する文書情報 9 0 2 に、当該文書の文書識別番号が存在するときは、その文書識別番号に対応する出現回数を 1 だけ加算する。また、登録しようとする単語が単語インデックス 2 0 5 のインデックス見出し語 9 0 1 に既に存在している場合において、インデックス見出し語 9 0 1 に対応する文書情報 9 0 2 に、当該文書の文書識別番号が存在しないときは、抽出テキストの文書識別番号を新たに登録し、出現回数を 1 とする。その後、ステップ S 4 0 0 8 に進む。

【0 0 7 7】

一方、前記ステップ S 4 0 0 6 の判別の結果、「未知語検索指定無し」である場合は、抽出された未知語を単語インデックス 2 0 5 に索引情報として登録することなく、ステップ S 4 0 0 8 へ進む。

【 0 0 7 8 】

ステップ S 4 0 0 8 では、単語抽出処理対象の文書（テキスト）のすべての文字列に対する処理を終えたか否かを、文書末尾に達したかどうかで判別する。その判別の結果、文書末尾に達していなければ前記ステップ S 4 0 0 1 へ戻って、未処理の文字列について上述の処理を繰り返す一方、文書末尾に達していれば本処理を終了する。

【 0 0 7 9 】

本実施の形態によれば、文字認識処理済み文書から抽出した未知語について、単語インデックス 2 0 5 に索引情報として登録するか否かを、未知語検索指定保持部 2 0 8 で指定できるようにしたので、検索漏れの抑制と検索精度の向上のいずれに重点を置くかという、検索者の意図を反映させて未知語の索引登録可否を決定することで、使い勝手がよく、検索漏れ及び誤検索の少ない適切な検索を行うことができる。

【 0 0 8 0 】

また、文書が文字認識処理済み文書でない場合、該文書から抽出された単語については、索引情報としての登録を一律に許可するようにしたので、高速な索引登録処理を実現することができる。

【 0 0 8 1 】

また、文字認識処理を経た文書について単語抽出処理を行う場合は、事前に文字校正処理を施すようにしたので、情報検索の精度を向上させることができる。

【 0 0 8 2 】

なお、本実施の形態では、文字認識処理済み文書から抽出した未知語については、「未知語検索指定有り」の場合にのみ、単語インデックス 2 0 5 に索引情報として登録されるようにしたが（図 1 4 のステップ S 4 0 0 6、S 4 0 0 7）、これに限るものでなく、文字認識処理済み文書から抽出した未知語については、一律に、単語インデックス 2 0 5 への索引情報としての登録を禁止するように構成してもよい。その場合は、図 1 4 のステップ S 4 0 0 6 の処理を省略すると共に、ステップ S 4 0 0 5 の判別の結果、単語抽出処理対象文書が文字認識結果文書（文字認識処理済み文書）である場合は、直ちに前記ステップ S 4 0 0 8 に進

むように処理すればよい。

【0083】

このようにすれば、単語抽出の対象となった文書の種別によって、単語インデックス 2 0 5 への索引情報としての登録の可否が決定されるので、例えば、文字認識処理済み文書から抽出した未知語については無駄で不適切な索引登録を抑制でき、索引登録処理時間の短縮化、索引サイズの縮小という効果が得られる。また、文字認識誤りの含まれ得る文書であっても、認識誤りに由来する不適切な索引登録を回避して、不適切な検索結果を抑制することができ、操作性に優れた、精度の高い情報検索が可能となる。一方では、文字認識誤りのない文書から抽出された単語については一律に登録可能とすることで、高速な索引登録処理を実現することができる。よって、文書の種別に応じて未知単語の索引登録の可否を決定することで、検索漏れ及び誤検索の少ない適切な検索を行うことができる。

【0084】

なお、本実施の形態では、単語インデックス 2 0 5、単語抽出辞書 2 0 7、文字校正用辞書 2 1 0 を単一の装置を構成する D I S K 1 4 に配置するものとして説明したが、これらの構成要素を異なる装置に分散配置し、N I C 1 9 を介してネットワーク上で処理を行うようにすることも可能である。

【0085】

なお、本発明は、複数の機器（例えばホストコンピュータ、インタフェース機器、リーダ、プリンタ等）から構成されるシステムに適用しても、ひとつの機器からなる装置（例えば、複写機、ファクシミリ装置等）に適用してもよい。

【0086】

また、本発明の目的は、実施の形態の機能を実現するソフトウェアのプログラムコードを記録した記憶媒体を、システム或いは装置に供給し、そのシステム或いは装置のコンピュータ（または C P U や M P U 等）が記憶媒体に格納されたプログラムコードを読み出して実行することによっても達成される。

【0087】

この場合、記憶媒体から読み出されたプログラムコード自体が前述した実施の形態の機能を実現することになり、そのプログラムコードを記憶した記憶媒体は

本発明を構成することになる。

【 0 0 8 8 】

又、プログラムコードを供給するための記憶媒体としては、例えば、フロッピー（登録商標）ディスク、ハードディスク、光磁気ディスク、CD-ROM、CD-R、CD-RW、DVD-ROM、DVD-RAM、DVD-RW、DVD+RW、磁気テープ、不揮発性のメモリカード、ROM等を用いることができる。

【 0 0 8 9 】

また、コンピュータが読み出したプログラムコードを実行することにより、上記実施の形態の機能が実現されるだけでなく、そのプログラムコードの指示に基づき、コンピュータ上で稼動しているOS（オペレーティングシステム）等が実際の処理の一部または全部を行い、その処理によって前述した実施の形態の機能が実現される場合も含まれる。

【 0 0 9 0 】

更に、記憶媒体から読み出されたプログラムコードが、コンピュータに挿入された機能拡張ボードやコンピュータに接続された機能拡張ユニットに備わるメモリに書き込まれた後、そのプログラムコードの指示に基づき、その機能拡張ボードや機能拡張ユニットに備わるCPU等が実際の処理の一部または全部を行い、その処理によって前述した実施の形態の機能が実現される場合も含まれる。

【 0 0 9 1 】

本発明の様々な例と実施例が示され説明されたが、当業者であれば、本発明の趣旨と範囲は本明細書の特定の説明と図に限定されるのではなく、本願特許請求の範囲にすべて述べられた様々の修正と変更にあふことが理解されるであろう。

【 0 0 9 2 】

本発明の実施態様の例を以下に列挙する。

【 0 0 9 3 】

〔実施態様1〕 文書から抽出された単語を前記文書と対応付けて文書検索用の索引情報として登録する索引情報登録手段と、該索引情報登録手段により登録された索引情報を参照して、要求された検索用情報に対応する文書を検索する文

書検索手段とを有する情報検索装置であって、前記文書から未知単語を抽出する未知単語抽出手段と、前記文書の種別を判別する文書種別判別手段と、前記文書種別判別手段により判別された文書の種別に応じて、前記未知単語抽出手段により抽出された未知単語の前記索引情報登録手段による索引情報としての登録の可否を決定する登録可否決定手段とを有することを特徴とする情報検索装置。

【0094】

〔実施態様2〕 前記文書種別判別手段は、前記文書が文字認識処理により得られた文字コードを含んだ文字認識処理済み文書であるか否かを判別し、前記登録可否決定手段は、前記文書が前記文字認識処理済み文書である場合は、該文書から前記未知単語抽出手段により抽出された未知単語の前記索引情報登録手段による索引情報としての登録を禁止することを特徴とする実施態様1記載の情報検索装置。

【0095】

〔実施態様3〕 前記文書種別判別手段は、前記文書が文字認識処理により得られた文字コードを含んだ文字認識処理済み文書であるか否かを判別するものであり、さらに、前記文書が前記文字認識処理済み文書である場合は、該文書から前記未知単語抽出手段により抽出された未知単語の前記索引情報登録手段による索引情報としての登録可否を指定する可否指定手段を有し、前記登録可否決定手段は、前記文書が前記文字認識処理済み文書である場合は、前記可否指定手段の指定に基づいて、前記未知単語抽出手段により抽出された未知単語の前記索引情報登録手段による索引情報としての登録の可否を決定することを特徴とする実施態様1記載の情報検索装置。

【0096】

〔実施態様4〕 前記文書種別判別手段は、前記文書が文字認識処理により得られた文字コードを含んだ文字認識処理済み文書であるか否かを判別し、前記登録可否決定手段は、前記文書が前記文字認識処理済み文書でない場合は、該文書から前記未知単語抽出手段により抽出された未知単語の前記索引情報登録手段による索引情報としての登録を許可する実施態様1～3のいずれかに記載の情報検索装置。

【0 0 9 7】

〔実施態様 5〕 前記文書種別判別手段は、前記文書が文字認識処理により得られた文字コードを含んだ文字認識処理済み文書であるか否かを判別するものであり、さらに、前記文書が前記文字認識処理済み文書である場合は、それに含まれている文字コードの誤りを校正する文字校正手段を有し、前記未知単語抽出手段は、前記文字校正手段により校正された後の文書から未知単語を抽出することを特徴とする実施態様 1 ～ 4 のいずれかに記載の情報検索装置。

【0 0 9 8】

〔実施態様 6〕 前記文書種別判別手段は、前記文書が画像文書であるか否かを判別するものであり、さらに、前記文書が画像文書である場合は、該画像文書から文字認識処理により文字コードを得る文字認識処理手段と、該文字認識処理手段により文字認識処理された文書に含まれている文字コードの誤りを校正する文字校正手段とを有し、前記未知単語抽出手段は、前記文字校正手段により校正された後の文書から未知単語を抽出することを特徴とする実施態様 1 ～ 4 のいずれかに記載の情報検索装置。

【0 0 9 9】

〔実施態様 7〕 前記文字校正手段は、類似の文字パターンを有する文字同士を対応付けて記憶した文字校正辞書と文字列照合のための単語辞書とを参照し、文字認識処理された文書に含まれている文字コードに前記文字校正辞書中の文字パターンを追加した文字候補列の中で、前記単語辞書に含まれる単語と一致する文字候補列を構成する文字に前記文字コードを変更することで、校正を行うことを特徴とする実施態様 5 または 6 記載の情報検索装置。

【0 1 0 0】

〔実施態様 8〕 前記文書種別判別手段は、前記文書のファイル名の拡張子及び前記文書に予め付与された属性情報の少なくとも一方に基づき文書の種別を判別することを特徴とする実施態様 1 ～ 7 のいずれかに記載の情報検索装置。

【0 1 0 1】

〔実施態様 9〕 文書から単語を抽出する単語抽出手段と、前記単語抽出手段により抽出された単語を前記文書と対応付けて、文書検索用の索引情報として登

録する索引情報登録手段と、前記索引情報登録手段により登録された索引情報を参照して、要求された検索用情報に対応する文書を検索する文書検索手段と、文字認識処理された文書に含まれている文字コードの誤りを校正する文字校正手段とを有し、前記単語抽出手段は、前記文字校正手段により校正された後の文書から単語を抽出することを特徴とする情報検索装置。

【0 1 0 2】

〔実施態様 1 0〕 前記文字校正手段は、類似の文字パターンを有する文字同士を対応付けて記憶した文字校正辞書と文字列照合のための単語辞書とを参照し、文字認識処理された文書に含まれている文字コードに前記文字校正辞書中の文字パターンを追加した文字候補列の中で、前記単語辞書に含まれる単語と一致する文字候補列を構成する文字に前記文字コードを変更することで、校正を行うことを特徴とする実施態様 9 記載の情報検索装置。

【0 1 0 3】

〔実施態様 1 1〕 文書から抽出された単語を前記文書と対応付けて文書検索用の索引情報として登録する索引情報登録ステップと、該索引情報登録ステップにより登録された索引情報を参照して、要求された検索用情報に対応する文書を検索する文書検索ステップとを有する情報検索方法であって、前記文書から未知単語を抽出する未知単語抽出ステップと、前記文書の種別を判別する文書種別判別ステップと、前記文書種別判別ステップにより判別された文書の種別に応じて、前記未知単語抽出ステップにより抽出された未知単語の前記索引情報登録ステップによる索引情報としての登録の可否を決定する登録可否決定ステップとを有することを特徴とする情報検索方法。

【0 1 0 4】

〔実施態様 1 2〕 文書から単語を抽出する単語抽出ステップと、前記単語抽出ステップにより抽出された単語を前記文書と対応付けて、文書検索用の索引情報として登録する索引情報登録ステップと、前記索引情報登録ステップにより登録された索引情報を参照して、要求された検索用情報に対応する文書を検索する文書検索ステップと、文字認識処理された文書に含まれている文字コードの誤りを校正する文字校正ステップとを有し、前記単語抽出ステップは、前記文字校正

ステップにより校正された後の文書から単語を抽出することを特徴とする情報検索方法。

【0 1 0 5】

〔実施態様 1 3〕 文書から抽出された単語を前記文書と対応付けて文書検索用の索引情報として登録する索引情報登録ステップと、該索引情報登録ステップにより登録された索引情報を参照して、要求された検索用情報に対応する文書を検索する文書検索ステップとをコンピュータに実行させる情報検索プログラムであって、前記文書から未知単語を抽出する未知単語抽出ステップと、前記文書の種別を判別する文書種別判別ステップと、前記文書種別判別ステップにより判別された文書の種別に応じて、前記未知単語抽出ステップにより抽出された未知単語の前記索引情報登録ステップによる索引情報としての登録の可否を決定する登録可否決定ステップとをコンピュータに実行させることを特徴とする情報検索プログラム。

【0 1 0 6】

〔実施態様 1 4〕 文書から単語を抽出する単語抽出ステップと、前記単語抽出ステップにより抽出された単語を前記文書と対応付けて、文書検索用の索引情報として登録する索引情報登録ステップと、前記索引情報登録ステップにより登録された索引情報を参照して、要求された検索用情報に対応する文書を検索する文書検索ステップと、文字認識処理された文書に含まれている文字コードの誤りを校正する文字校正ステップとをコンピュータに実行させる情報検索プログラムであって、前記単語抽出ステップは、前記文字校正ステップにより校正された後の文書から単語を抽出することを特徴とする情報検索プログラム。

【0 1 0 7】

〔実施態様 1 5〕 文書から抽出された単語を前記文書と対応付けて文書検索用の索引情報として登録する索引情報登録ステップと、該索引情報登録ステップにより登録された索引情報を参照して、要求された検索用情報に対応する文書を検索する文書検索ステップとをコンピュータに実行させる情報検索プログラムを記憶したコンピュータ読み取り可能な記憶媒体であって、前記文書から未知単語を抽出する未知単語抽出ステップと、前記文書の種別を判別する文書種別判別ス

テップと、前記文書種別判別ステップにより判別された文書の種別に応じて、前記未知単語抽出ステップにより抽出された未知単語の前記索引情報登録ステップによる索引情報としての登録の可否を決定する登録可否決定ステップとをコンピュータに実行させる情報検索プログラムを記憶したことを特徴とする記憶媒体。

【0 1 0 8】

〔実施態様 1 6〕 文書から単語を抽出する単語抽出ステップと、前記単語抽出ステップにより抽出された単語を前記文書と対応付けて、文書検索用の索引情報として登録する索引情報登録ステップと、前記索引情報登録ステップにより登録された索引情報を参照して、要求された検索用情報に対応する文書を検索する文書検索ステップと、文字認識処理された文書に含まれている文字コードの誤りを校正する文字校正ステップとをコンピュータに実行させる情報検索プログラムを記憶し、前記単語抽出ステップは、前記文字校正ステップにより校正された後の文書から単語を抽出することを特徴とする記憶媒体。

【0 1 0 9】

【発明の効果】

以上説明したように、本発明の請求項 1 によれば、文書の種別に応じて未知単語の索引登録の可否を決定することで、検索漏れ及び誤検索の少ない適切な検索を行うことができる。

【0 1 1 0】

本発明の請求項 2 によれば、校正後に単語抽出することで、検索精度を向上させることができる。

【図面の簡単な説明】

【図 1】

本発明の一実施の形態に係る情報検索装置の全体構成を示すブロック図である。

【図 2】

本発明の実施の形態の情報検索装置の機能構成を示すブロック図である。

【図 3】

検索語を指示する場合において、表示器 C R T に画面表示される操作パネルの

例を示す図である。

【図 4】

検索要求入力処理部による検索処理実行後の検索結果の表示例を示す図である。

。

【図 5】

文字を表す画像に対して文字認識処理を施した結果得られた抽出テキストを示す図である。

【図 6】

図 5 に例示された文字認識誤りを含んでいる抽出テキストに対して文字校正処理を行った場合の結果を示す図である。

【図 7】

文字校正用辞書の構成を示す概念図である。

【図 8】

単語抽出辞書の構成を示す概念図である。

【図 9】

単語インデックスの構成を示す概念図である。

【図 1 0】

文字校正処理で用いられる文字候補ラティスの構成を示す概念図である。

【図 1 1】

本実施の形態における文書検索処理のフローチャートを示す図である。

【図 1 2】

本実施の形態における文書登録処理のフローチャートを示す図である。

【図 1 3】

図 1 2 のステップ S 3 0 0 3 で実行される文字校正処理のフローチャートを示す図である。

【図 1 4】

図 1 2 のステップ S 3 0 0 4 で実行される単語抽出処理のフローチャートを示す図である。

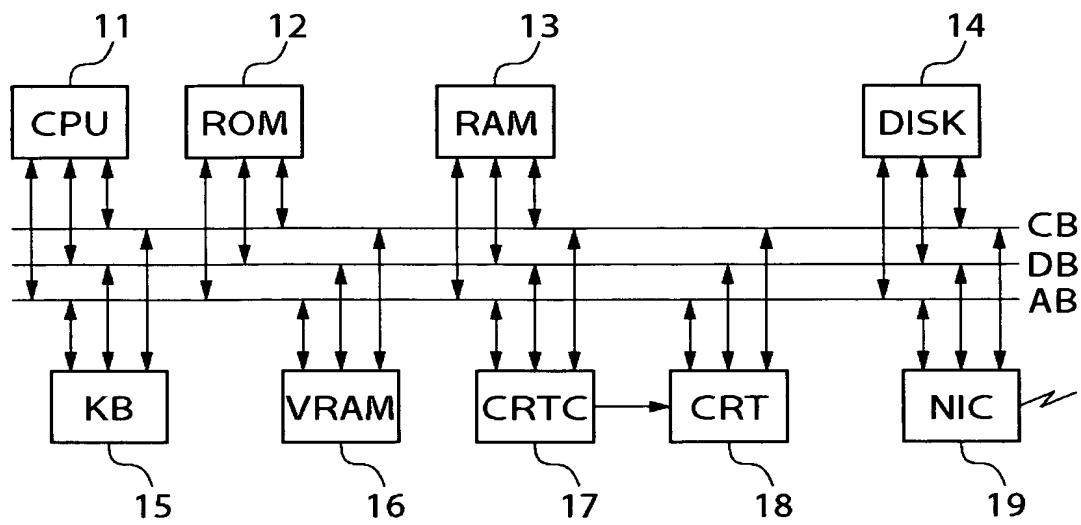
【符号の説明】

1 1. C P U (索引情報登録手段、文書検索手段、未知単語抽出手段、単語抽出手段、文書種別判別手段、文字認識処理手段、登録可否決定手段、文字校正手段)

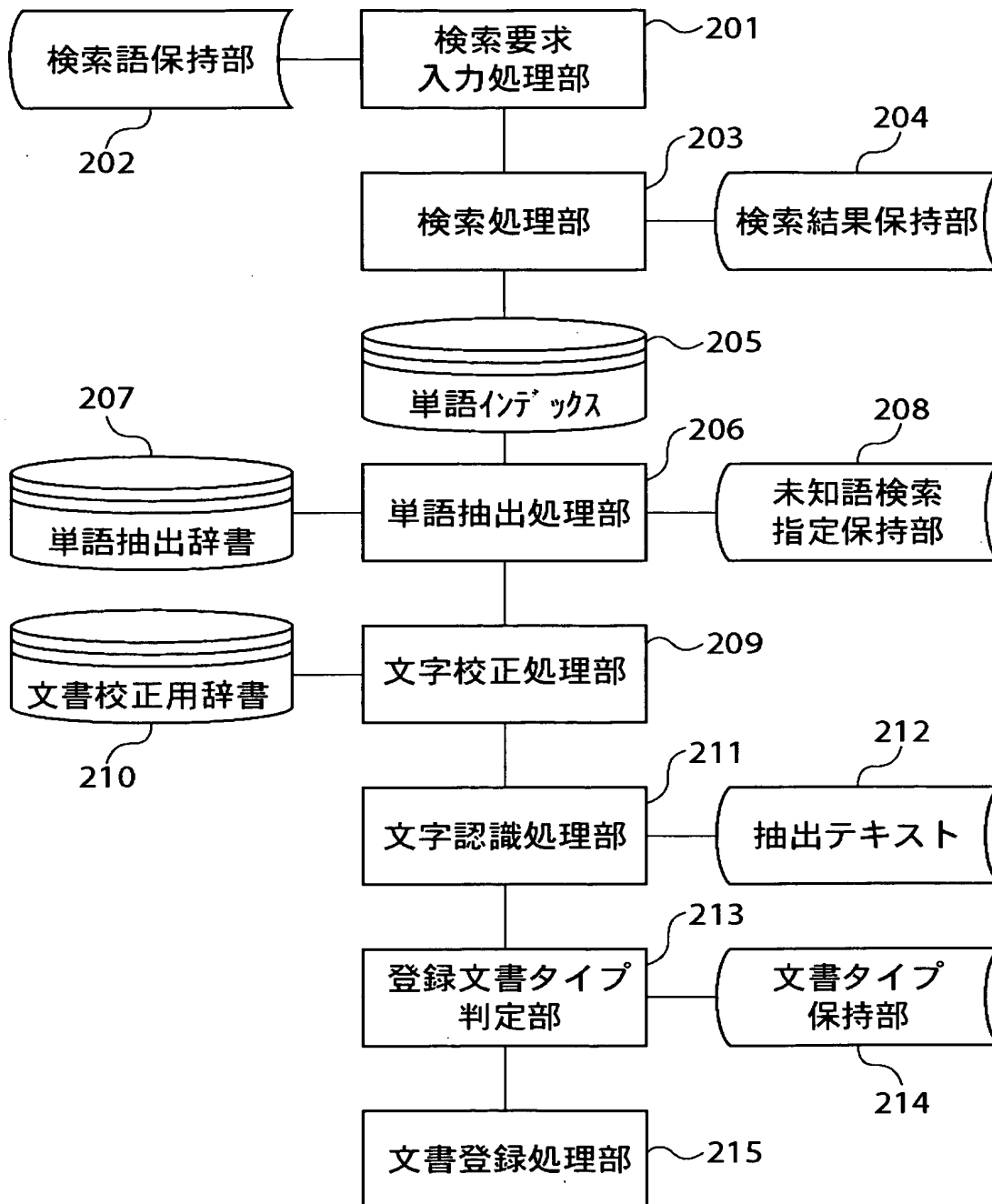
- 2 0 1 検索要求入力処理部
- 2 0 2 検索語保持部
- 2 0 3 検索処理部 (文書検索手段)
- 2 0 4 検索結果保持部
- 2 0 5 単語インデックス
- 2 0 6 単語抽出処理部 (未知単語抽出手段、単語抽出手段)
- 2 0 7 単語抽出辞書 (単語辞書)
- 2 0 8 未知語検索指定保持部 (可否指定手段)
- 2 0 9 文字校正処理部 (文字校正手段)
- 2 1 0 文字校正用辞書
- 2 1 1 文字認識処理部 (文字認識処理手段)
- 2 1 2 抽出テキスト
- 2 1 3 登録文書タイプ判定部 (文書種別判別手段)
- 2 1 4 文書タイプ保持部
- 2 1 5 文書登録処理部

【書類名】 図面

【図 1】



【図 2】



【図 3】

類似文書検索

301

302

検索文入力

カラープリンタの売上

304(カ-リル)

303

308 検索

309 キャンセル

310

【図 4】

301

302

類似文書検索

検索文入力

303

カラープリンタの売上↑

304

402

403

404

310

405

308

検索

309

キャンセル

406

ランク	文書表題	ファイル名	内容
25	国内市場占有率調査結果	ANNOUNCE.HTML	2001年度複写機売上▲
26	国内プリンタ販売動向報告書	1997REPORT.txt	1997年国内プリンタ販売動向報告書1998
27	マーケティング概論	MARKETING.doc	最新のマーケティング理論を学び実践において
28	カラーマネージメント	TQC5213.HTML	複数の機器間での色合わせをすることが重要となる
29	カラープリンタマニュアル	MP970621.pdf	さまざまな遠隔サービスを提供いたしますため
30	インターネットマーケティングビジネスの実際	WC990123.doc	LANケーブルをはいまわすことなく、高速ネッ▼

【図 5】

.....
このモノクロの写真に写っているモルルールは当時の重要な交通手段であった.....

【図 6】

.....
このモノクロの写真に写っているモルルールは当時の重要な交通手段であった.....

【図 7】

カカカ刀

口口
へへ
エエ
二二
.....

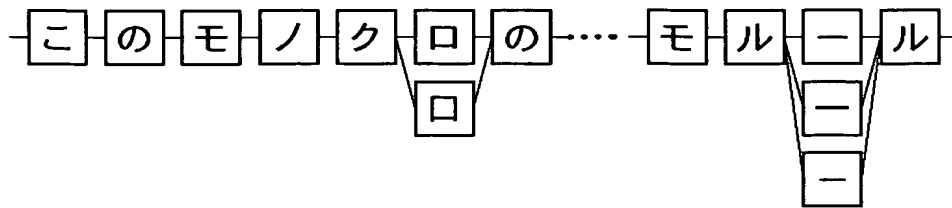
【図 8】

見出し語	品詞
...	...
カラー	名詞
...	...
モノクロ	名詞
...	...

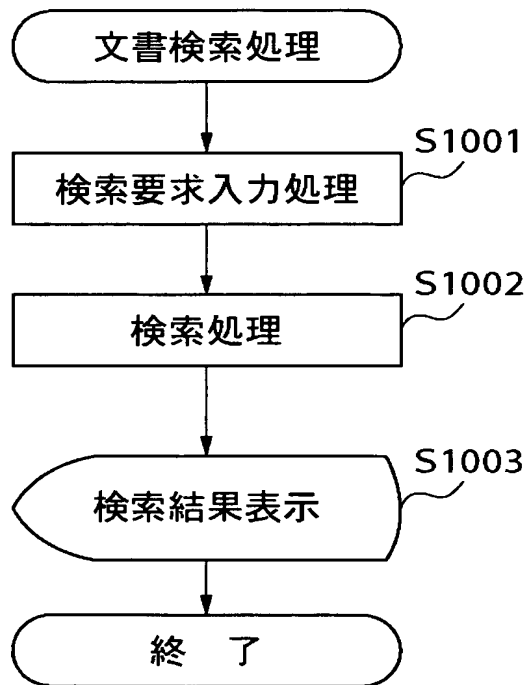
【図 9】

インデックス見出し語	文書情報
...	...
カラー	(1000,15),(1200,5),(1252,8),(1800,3)
カレンダー	(328,2),(512,1),(999,2)
...	...
モノクロ	(2100,1)
...	...

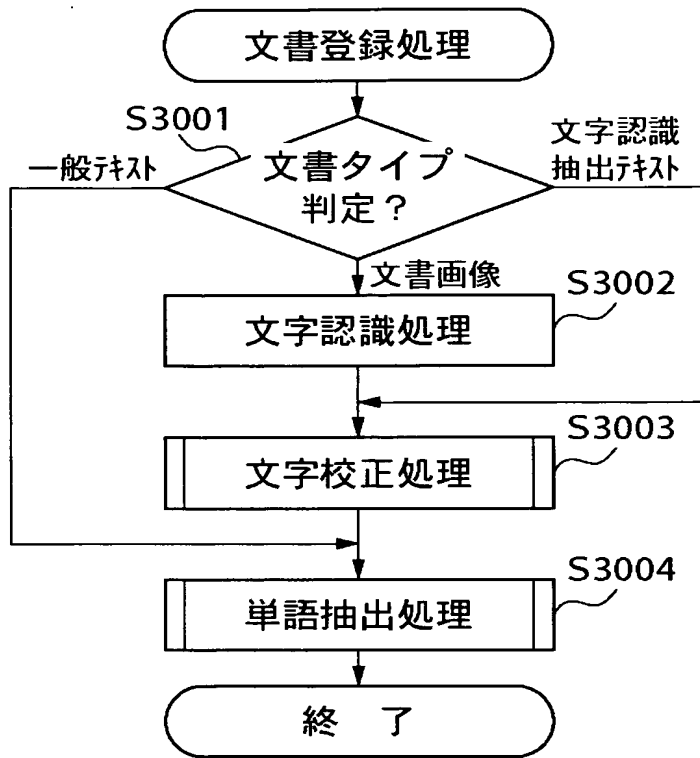
【図 10】



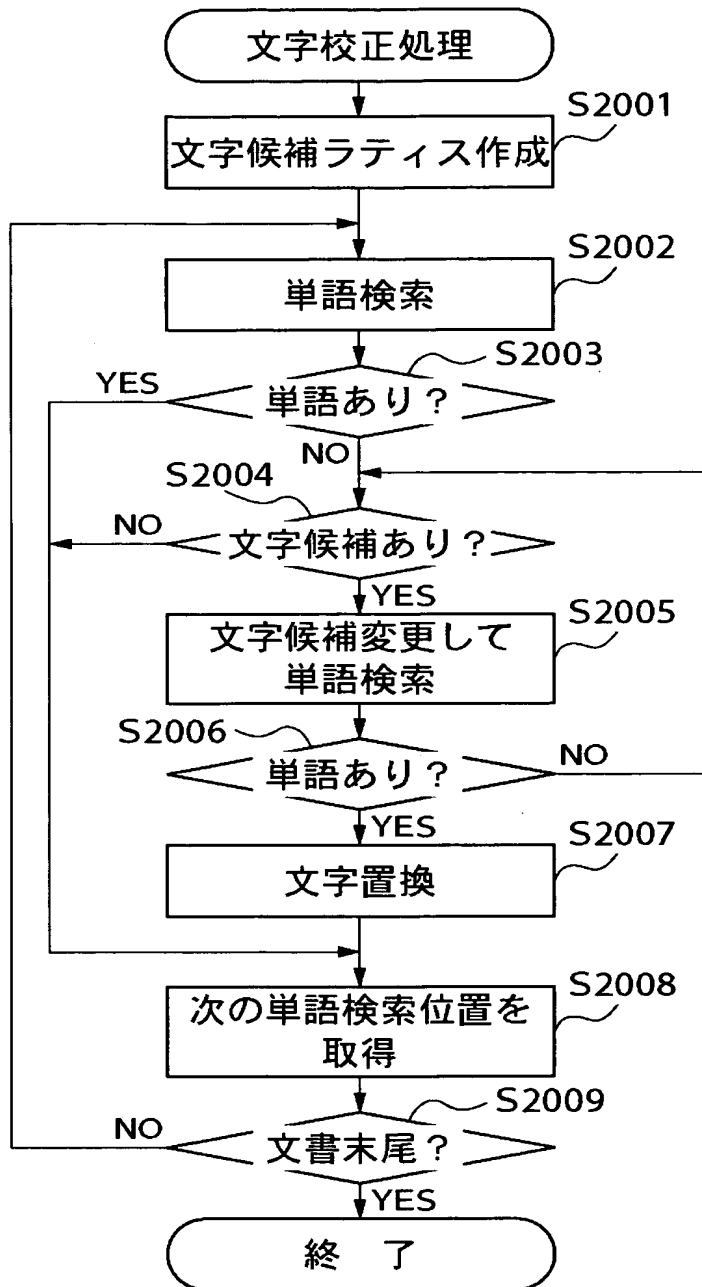
【図 11】



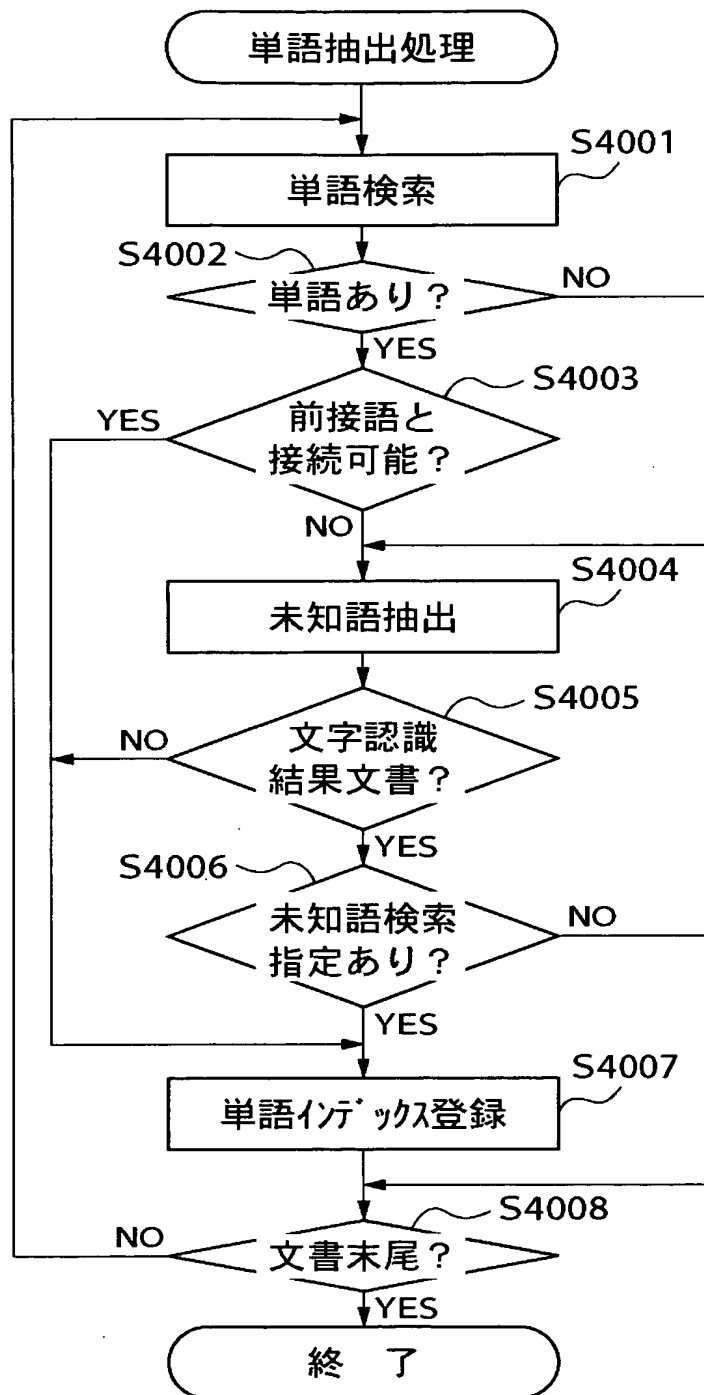
【図 12】



【図 13】



【図 14】



【書類名】 要約書

【要約】

【課題】 検索漏れ及び誤検索の少ない適切な検索を行えるようにする。

【解決手段】 単語抽出処理対象の文書から単語検索、または辞書登録されていない未知語の抽出を行い、単語抽出処理対象文書が文字認識結果文書、すなわち、文字認識処理済み文書でない場合は、単語インデックス 2 0 5 に索引情報として登録する。一方、文字認識処理済み文書である場合は、未知語検索指定保持部 2 0 8 を参照することで、未知語検索指定の有無を判定する。そして、「未知語検索指定有り」である場合は、抽出された未知語を、単語インデックス 2 0 5 に索引情報として登録する一方、「未知語検索指定無し」である場合は、索引情報として登録しない。

【選択図】 図 1 4

特願 2 0 0 3 - 0 1 3 4 2 8

出 願 人 履 歴 情 報

識別番号

[0 0 0 0 0 1 0 0 7]

1. 変更年月日

1 9 9 0 年 8 月 3 0 日

[変更理由]

新規登録

住 所

東京都大田区下丸子 3 丁目 3 0 番 2 号

氏 名

キヤノン株式会社